



**Manuscript version: Author's Accepted Manuscript**

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

**Persistent WRAP URL:**

<http://wrap.warwick.ac.uk/137739>

**How to cite:**

Please refer to published version for the most recent bibliographic citation information. If a published version is known of, the repository item page linked to above, will contain details on accessing it.

**Copyright and reuse:**

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

**Publisher's statement:**

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: [wrap@warwick.ac.uk](mailto:wrap@warwick.ac.uk).

# Bootstrap-Based Inference for Cube Root Asymptotics\*

Matias D. Cattaneo<sup>†</sup>

Michael Jansson<sup>‡</sup>

Kenichi Nagasawa<sup>§</sup>

May 29, 2020

## Abstract

This paper proposes a valid bootstrap-based distributional approximation for  $M$ -estimators exhibiting a Chernoff (1964)-type limiting distribution. For estimators of this kind, the standard nonparametric bootstrap is inconsistent. The method proposed herein is based on the nonparametric bootstrap, but restores consistency by altering the shape of the criterion function defining the estimator whose distribution we seek to approximate. This modification leads to a generic and easy-to-implement resampling method for inference that is conceptually distinct from other available distributional approximations. We illustrate the applicability of our results with four examples in econometrics and machine learning.

*Keywords:* cube root asymptotics, bootstrapping, maximum score, empirical risk minimization.

---

\*A previous version of this paper circulated under the title “Bootstrap-Based Inference for Cube Root Consistent Estimators”. For comments and suggestions, we are grateful to Mehmet Caner, Andreas Hagemann, Kei Hirano, Bo Honoré, Guido Imbens, Guido Kuersteiner, Mykhaylo Shkolnikov, Ronnie Sircar, Ulrich Müller, Whitney Newey, and participants at various conferences, workshops, and seminars. Cattaneo gratefully acknowledges financial support from the National Science Foundation through grants SES-1459931 and SES-1947805, and Jansson gratefully acknowledges financial support from the National Science Foundation through grants SES-1459967 and SES-1947662 and the research support of CREATES (funded by the Danish National Research Foundation under grant no. DNRF78).

<sup>†</sup>Department of Operations Research and Financial Engineering, Princeton University.

<sup>‡</sup>Department of Economics, University of California at Berkeley and CREATES.

<sup>§</sup>Department of Economics, University of Warwick.

# 1 Introduction

In a seminal paper, [Kim and Pollard \(1990\)](#) studied estimators exhibiting “cube root asymptotics”. These estimators not only have a non-standard rate of convergence, but also have the property that rather than being Gaussian their limiting distributions are of [Chernoff \(1964\)](#) type; i.e., the non-Gaussian limiting distribution is that of the maximizer of a Gaussian process. Kim and Pollard’s results cover not only celebrated examples such as maximum score estimator of [Manski \(1975\)](#) and the isotonic density estimator of [Grenander \(1956\)](#), but also more contemporary estimators arising in examples related to classification problems in machine learning ([Mohammadi and van de Geer, 2005](#)), nonparametric inference under shape restrictions ([Groeneboom and Jongbloed, 2018](#)), massive data  $M$ -estimation framework ([Shi, Lu, and Song, 2018](#)), and maximum score estimation in high-dimensional settings ([Mukherjee, Banerjee, and Ritov, 2019](#)). Moreover, [Seo and Otsu \(2018\)](#) recently generalized [Kim and Pollard \(1990\)](#) to allow for  $n$ -varying objective functions ( $n$  denotes the sample size), further widening the applicability of cube-root-type asymptotics. For example, their results cover the conditional maximum score estimator of [Honoré and Kyriazidou \(2000\)](#).

An important feature of Chernoff-type asymptotic distributional approximations is that the covariance kernel of the Gaussian process characterizing the limiting distribution often depends on an infinite-dimensional nuisance parameter. From the perspective of inference, this feature of the limiting distribution represents a nontrivial complication relative to the conventional asymptotically normal case, where the limiting distribution is known up to the value of a finite-dimensional nuisance parameter (namely, the covariance matrix of the limiting distribution). The dependence of the limiting distribution on an infinite-dimensional nuisance parameter implies that resampling-based distributional approximations seem to offer the most attractive approach to inference in estimation problems exhibiting cube root asymptotics. Unfortunately, however, the standard nonparametric bootstrap is well known to be invalid in this setting ([Abrevaya and Huang, 2005](#); [Léger and MacGibbon, 2006](#); [Kosorok, 2008](#); [Sen, Banerjee, and Woodroffe, 2010](#)). The purpose of this paper is to propose a generic and easy-to-implement bootstrap-based distributional approximation applicable in the context of cube root asymptotics.

As does the familiar nonparametric bootstrap, the method proposed herein employs bootstrap samples of size  $n$  from the empirical distribution function. But unlike the nonparametric bootstrap,

which is inconsistent, our method offers a consistent distributional approximation for estimators exhibiting cube root asymptotics. Consistency is achieved by altering the shape of the criterion function defining the estimator whose distribution we seek to approximate. Heuristically, the method is designed to ensure that the bootstrap version of a certain empirical process has a mean resembling the large sample version of its population counterpart. The latter is quadratic in the problems we study, and known up to the value of a certain matrix. As a consequence, the only ingredient needed to implement the proposed “reshapement” of the objective function is a consistent estimator of the unknown matrix entering the quadratic mean of the empirical process. Such estimators turn out to be generically available and easy to compute.

This paper is not the first to propose a consistent resampling-based distributional approximation for cube-root-type estimators. For canonical cube root asymptotic problems, the best known consistent alternative to the nonparametric bootstrap is probably subsampling ([Politis and Romano, 1994](#)), whose applicability was pointed out by [Delgado, Rodriguez-Poo, and Wolf \(2001\)](#). Related applicable methods are the  $m$  out of  $n$  bootstrap ([Bickel, Götze, and van Zwet, 1997](#)), whose applicability was discussed and extended by [Lee and Pun \(2006\)](#) and [Lee and Yang \(2020\)](#), the rescaled bootstrap ([Dümbgen, 1993](#)), and the numerical bootstrap ([Hong and Li, 2020](#)). In addition, case-specific (smooth or non-standard) bootstrap methods have been proposed for leading examples such as monotone density estimation ([Kosorok, 2008](#); [Sen, Banerjee, and Woodroffe, 2010](#)), maximum score estimation ([Patra, Seijo, and Sen, 2018](#)), and the current status model ([Groeneboom and Hendrickx, 2018](#)). For the more generic cube-root-type estimators analyzed in [Seo and Otsu \(2018\)](#), subsampling appears to be the only method available, and indeed the authors discuss in their concluding remarks the need for (and importance of) developing resampling methods based on the standard nonparametric bootstrap. Our paper appears to be the first to provide one such method.

Like ours, each of the resampling methods mentioned above can be viewed as offering a “robust” alternative to the standard nonparametric bootstrap but, unlike ours, existing methods achieve consistency by modifying the distribution used to generate the bootstrap sample. In contrast, our bootstrap-based method achieves consistency by means of an analytic modification of the objective function used to construct the bootstrap-based distributional approximation. As further discussed below, this approach results in a procedure that is conceptually related to the bootstrap

methods developed by [Andrews and Soares \(2010\)](#) and [Fang and Santos \(2019\)](#) in other econometrics contexts.

Implementation of our procedure is not computationally demanding. Indeed, the only ingredient needed to implement our modification on the objective function is a consistent estimator of a certain Hessian matrix. We propose a generic estimator based on numerical derivatives and present a consistency result as well as an approximate mean square error expansion for that estimator. In addition, we illustrate how example-specific features can be sometimes exploited to construct alternative estimators.

The paper proceeds as follows. Section 2 is heuristic and outlines the main idea underlying our approach in the  $M$ -estimation setting of [Kim and Pollard \(1990\)](#). Section 3 then makes that heuristic discussion rigorous in a more general setting similar to that of [Seo and Otsu \(2018\)](#). Section 4 illustrates our bootstrap-based inference method with four examples: the maximum score estimator of [Manski \(1975, 1985\)](#), the conditional maximum score panel data estimator of [Manski \(1987\)](#), the conditional maximum score dynamic panel data estimator of [Honoré and Kyriazidou \(2000\)](#), and the classification estimator of [Mohammadi and van de Geer \(2005\)](#). Section 5 reports simulation evidence for the case of the maximum score estimator, and Section 6 concludes. Section 7 describes the proof of our main result, while the supplemental appendix contains omitted proofs and details.

## 2 Heuristics

Suppose  $\theta_0 \in \Theta \subseteq \mathbb{R}^d$  is an estimand admitting the characterization

$$\theta_0 = \operatorname{argmax}_{\theta \in \Theta} M_0(\theta), \quad M_0(\theta) = \mathbb{E}[m_0(\mathbf{z}, \theta)], \quad (1)$$

where  $m_0$  is a known function, and where  $\mathbf{z}$  is a random vector of which a random sample  $\mathbf{z}_1, \dots, \mathbf{z}_n$  is available. Studying estimation problems of this kind for non-smooth  $m_0$ , [Kim and Pollard \(1990\)](#) gave conditions under which the  $M$ -estimator

$$\hat{\theta}_n = \operatorname{argmax}_{\theta \in \Theta} \hat{M}_n(\theta), \quad \hat{M}_n(\theta) = \frac{1}{n} \sum_{i=1}^n m_0(\mathbf{z}_i, \theta),$$

exhibits cube root asymptotics:

$$\sqrt[3]{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \rightsquigarrow \operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} \{\mathcal{G}_0(\mathbf{s}) + \mathcal{Q}_0(\mathbf{s})\}, \quad (2)$$

where  $\rightsquigarrow$  denotes weak convergence,  $\mathcal{G}_0$  is a non-degenerate zero-mean Gaussian process, and  $\mathcal{Q}_0(\mathbf{s}) = -\mathbf{s}'\mathbf{H}_0\mathbf{s}/2$ , where  $\mathbf{H}_0 = -\partial^2 M_0(\boldsymbol{\theta}_0)/\partial\boldsymbol{\theta}\partial\boldsymbol{\theta}'$ .

Whereas the matrix  $\mathbf{H}_0$  governing the shape of  $\mathcal{Q}_0$  is finite-dimensional, the covariance kernel of  $\mathcal{G}_0$  in (2) typically involves infinite-dimensional unknown quantities. As a consequence, the limiting distribution of  $\hat{\boldsymbol{\theta}}_n$  tends to be more difficult to approximate than Gaussian distributions, implying in turn that basing inference on  $\hat{\boldsymbol{\theta}}_n$  is more challenging under cube root asymptotics than in the more familiar case where  $\hat{\boldsymbol{\theta}}_n$  is  $\sqrt{n}$ -consistent and asymptotically normally distributed.

As a candidate method of approximating the distribution of  $\hat{\boldsymbol{\theta}}_n$ , consider the nonparametric bootstrap. To describe it, let  $\mathbf{z}_{1,n}^*, \dots, \mathbf{z}_{n,n}^*$  denote a random sample from the empirical distribution of  $\mathbf{z}_1, \dots, \mathbf{z}_n$  and let the natural bootstrap analogue of  $\hat{\boldsymbol{\theta}}_n$  be denoted by

$$\hat{\boldsymbol{\theta}}_n^* = \operatorname{argmax}_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \hat{M}_n^*(\boldsymbol{\theta}), \quad \hat{M}_n^*(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n m_0(\mathbf{z}_{i,n}^*, \boldsymbol{\theta}).$$

Then, the nonparametric bootstrap estimator of  $\mathbb{P}[\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0 \leq \cdot]$  is given by  $\mathbb{P}_n^*[\hat{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n \leq \cdot]$ , where  $\mathbb{P}_n^*$  denotes a probability computed under the bootstrap distribution conditional on the data. As is well documented, however, this estimator is inconsistent under cube root asymptotics (Abrevaya and Huang, 2005; Léger and MacGibbon, 2006; Kosorok, 2008; Sen, Banerjee, and Woodroffe, 2010).

For the purpose of giving a heuristic, yet constructive, explanation of the inconsistency of the nonparametric bootstrap, it is helpful to recall that a proof of (2) can be based on the representation

$$\sqrt[3]{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = \operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} \{\hat{G}_n(\mathbf{s}) + Q_n(\mathbf{s})\}, \quad (3)$$

where, for  $\mathbf{s}$  such that  $\boldsymbol{\theta}_0 + \mathbf{s}n^{-1/3} \in \boldsymbol{\Theta}$ ,

$$\hat{G}_n(\mathbf{s}) = n^{2/3}[\hat{M}_n(\boldsymbol{\theta}_0 + \mathbf{s}n^{-1/3}) - \hat{M}_n(\boldsymbol{\theta}_0) - M_0(\boldsymbol{\theta}_0 + \mathbf{s}n^{-1/3}) + M_0(\boldsymbol{\theta}_0)] \quad (4)$$

is a zero-mean random process, while

$$Q_n(\mathbf{s}) = n^{2/3}[M_0(\boldsymbol{\theta}_0 + \mathbf{s}n^{-1/3}) - M_0(\boldsymbol{\theta}_0)] \quad (5)$$

is a non-random function that is correctly centered in the sense that  $\operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} Q_n(\mathbf{s}) = \mathbf{0}$ . In cases where  $m_0$  is non-smooth but  $M_0$  is smooth,  $\hat{G}_n$  and  $Q_n$  are usually asymptotically Gaussian

and asymptotically quadratic, respectively, in the sense that

$$\hat{G}_n(\mathbf{s}) \rightsquigarrow \mathcal{G}_0(\mathbf{s}) \quad (6)$$

and

$$Q_n(\mathbf{s}) \rightarrow \mathcal{Q}_0(\mathbf{s}). \quad (7)$$

Under regularity conditions ensuring among other things that the convergence in (6) and (7) is suitably uniform in  $\mathbf{s}$ , (2) then follows from an application of a continuous mapping-type theorem for the argmax functional to the representation in (3).

Similarly to (3), the bootstrap analogue of  $\hat{\boldsymbol{\theta}}_n$  admits a representation of the form

$$\sqrt[3]{n}(\hat{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n) = \operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} \{\hat{G}_n^*(\mathbf{s}) + \hat{Q}_n(\mathbf{s})\},$$

where, for  $\mathbf{s}$  such that  $\hat{\boldsymbol{\theta}}_n + \mathbf{s}n^{-1/3} \in \boldsymbol{\Theta}$ ,

$$\hat{G}_n^*(\mathbf{s}) = n^{2/3}[\hat{M}_n^*(\hat{\boldsymbol{\theta}}_n + \mathbf{s}n^{-1/3}) - \hat{M}_n^*(\hat{\boldsymbol{\theta}}_n) - \hat{M}_n(\hat{\boldsymbol{\theta}}_n + \mathbf{s}n^{-1/3}) + \hat{M}_n(\hat{\boldsymbol{\theta}}_n)]$$

and

$$\hat{Q}_n(\mathbf{s}) = n^{2/3}[\hat{M}_n(\hat{\boldsymbol{\theta}}_n + \mathbf{s}n^{-1/3}) - \hat{M}_n(\hat{\boldsymbol{\theta}}_n)].$$

Under mild conditions,  $\hat{G}_n^*$  satisfies the following bootstrap counterpart of (6):

$$\hat{G}_n^*(\mathbf{s}) \rightsquigarrow_{\mathbb{P}} \mathcal{G}_0(\mathbf{s}), \quad (8)$$

where  $\rightsquigarrow_{\mathbb{P}}$  denotes conditional weak convergence in probability (defined as in [van der Vaart and Wellner, 1996](#), Section 2.9). On the other hand, although  $\hat{Q}_n$  is non-random under the bootstrap distribution and satisfies  $\operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} \hat{Q}_n(\mathbf{s}) = \mathbf{0}$ , it turns out that  $\hat{Q}_n(\mathbf{s}) \not\rightarrow_{\mathbb{P}} \mathcal{Q}_0(\mathbf{s})$  in general. In other words, the natural bootstrap counterpart of (7) typically fails and, as a partial consequence, so does the natural bootstrap counterpart of (2); i.e.,  $\sqrt[3]{n}(\hat{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n) \not\rightsquigarrow_{\mathbb{P}} \operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} \{\mathcal{G}_0(\mathbf{s}) + \mathcal{Q}_0(\mathbf{s})\}$ .

To the extent that the inconsistency of the bootstrap can be attributed to the fact that the shape of  $\hat{Q}_n$  fails to replicate that of  $Q_n$ , it seems plausible that a consistent bootstrap-based distributional approximation can be obtained by basing the approximation on

$$\tilde{\boldsymbol{\theta}}_n^* = \operatorname{argmax}_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \tilde{M}_n^*(\boldsymbol{\theta}), \quad \tilde{M}_n^*(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \tilde{m}_n(\mathbf{z}_{i,n}^*, \boldsymbol{\theta}),$$

where  $\tilde{m}_n$  is a suitably “reshaped” version of  $m_0$  satisfying two properties. First,  $\tilde{G}_n^*$  should be

asymptotically equivalent to  $\hat{G}_n^*$ , where  $\tilde{G}_n^*$  is the counterpart of  $\hat{G}_n^*$  associated with  $\tilde{m}_n$  :

$$\tilde{G}_n^*(\mathbf{s}) = n^{2/3}[\tilde{M}_n^*(\hat{\boldsymbol{\theta}}_n + \mathbf{s}n^{-1/3}) - \tilde{M}_n^*(\hat{\boldsymbol{\theta}}_n) - \tilde{M}_n(\hat{\boldsymbol{\theta}}_n + \mathbf{s}n^{-1/3}) + \tilde{M}_n(\hat{\boldsymbol{\theta}}_n)], \quad \tilde{M}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \tilde{m}_n(\mathbf{z}_i, \boldsymbol{\theta}).$$

Second, and most importantly,  $\tilde{Q}_n$  should be asymptotically quadratic, where  $\tilde{Q}_n$  is the counterpart of  $\hat{Q}_n$  associated with  $\tilde{m}_n$ :

$$\tilde{Q}_n(\mathbf{s}) = n^{2/3}[\tilde{M}_n(\hat{\boldsymbol{\theta}}_n + \mathbf{s}n^{-1/3}) - \tilde{M}_n(\hat{\boldsymbol{\theta}}_n)].$$

Accordingly, let

$$\tilde{m}_n(\mathbf{z}, \boldsymbol{\theta}) = m_0(\mathbf{z}, \boldsymbol{\theta}) - \hat{M}_n(\boldsymbol{\theta}) - \frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n)' \tilde{\mathbf{H}}_n(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n),$$

where  $\tilde{\mathbf{H}}_n$  is an estimator of  $\mathbf{H}_0$ . Then

$$\sqrt[3]{n}(\tilde{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n) = \operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} \{\tilde{G}_n^*(\mathbf{s}) + \tilde{Q}_n(\mathbf{s})\},$$

where, by construction,  $\tilde{G}_n^*(\mathbf{s}) = \hat{G}_n^*(\mathbf{s})$  and  $\tilde{Q}_n(\mathbf{s}) = -\mathbf{s}' \tilde{\mathbf{H}}_n \mathbf{s} / 2$ . Because  $\tilde{G}_n^* = \hat{G}_n^*$ ,  $\tilde{G}_n^*(\mathbf{s}) \rightsquigarrow_{\mathbb{P}} \mathcal{G}_0(\mathbf{s})$  whenever (8) holds. In addition,  $\tilde{Q}_n(\mathbf{s}) \rightarrow_{\mathbb{P}} \mathcal{Q}_0(\mathbf{s})$  provided  $\tilde{\mathbf{H}}_n \rightarrow_{\mathbb{P}} \mathbf{H}_0$ . As a consequence, it seems plausible that if  $\tilde{\mathbf{H}}_n \rightarrow_{\mathbb{P}} \mathbf{H}_0$ , then  $\sqrt[3]{n}(\tilde{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n) \rightsquigarrow_{\mathbb{P}} \operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} \{\mathcal{G}_0(\mathbf{s}) + \mathcal{Q}_0(\mathbf{s})\}$ .

For the purposes of situating this paper in the literature, the following alternative heuristic explanation of our approach may be useful. Restating the result in (2) as

$$\sqrt[3]{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \rightsquigarrow \mathcal{S}_0(\mathcal{G}_0), \quad \mathcal{S}_0(\mathcal{G}) = \operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} \{\mathcal{G}(\mathbf{s}) + \mathcal{Q}_0(\mathbf{s})\},$$

our procedure approximates the distribution of  $\mathcal{S}_0(\mathcal{G}_0)$  by that of  $\tilde{\mathcal{S}}_n(\hat{G}_n^*)$ , where the distribution of the bootstrap process  $\hat{G}_n^*$  approximates that of  $\mathcal{G}_0$  and where  $\tilde{\mathcal{S}}_n(\mathcal{G}) = \operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} \{\mathcal{G}(\mathbf{s}) + \tilde{Q}_n(\mathbf{s})\}$  is an estimator of  $\mathcal{S}_0(\mathcal{G})$ . In other words, our procedure replaces the functional  $\mathcal{S}_0$  with a consistent estimator (namely,  $\tilde{\mathcal{S}}_n$ ) and its random argument  $\mathcal{G}_0$  with a bootstrap approximation (namely,  $\hat{G}_n^*$ ). The same type of generic construction has appeared in the econometrics literature before, notably in [Andrews and Soares \(2010\)](#) and [Fang and Santos \(2019\)](#).

Our bootstrap-based distributional approximation can be shown to be consistent also in the more standard case where  $m_n(\mathbf{z}, \boldsymbol{\theta})$  is sufficiently smooth in  $\boldsymbol{\theta}$  to ensure that an approximate maximizer of  $\hat{M}_n$  is asymptotically normal and that the nonparametric bootstrap is consistent. In fact,  $\tilde{\boldsymbol{\theta}}_n^*$  is (first-order) asymptotically equivalent to  $\hat{\boldsymbol{\theta}}_n^*$  in that standard case, so our procedure can be interpreted as a modification of the nonparametric bootstrap that is designed to be “robust” to



the types of non-smoothness that give rise to cube root asymptotics.

### 3 Main Result

When making the heuristics of Section 2 precise, we consider the more general situation where the estimator  $\hat{\boldsymbol{\theta}}_n$  is an approximate maximizer (with respect to  $\boldsymbol{\theta} \in \boldsymbol{\Theta} \subseteq \mathbb{R}^d$ ) of

$$\hat{M}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n m_n(\mathbf{z}_i, \boldsymbol{\theta}),$$

where  $m_n$  is a known function, and where  $\mathbf{z}_1, \dots, \mathbf{z}_n$  is a random sample of a random vector  $\mathbf{z}$ . This formulation of  $\hat{M}_n$ , which reduces to that of Section 2 when  $m_n$  does not depend on  $n$ , is adopted in order to cover certain estimation problems where, rather than admitting a characterization of the form (1), the estimand  $\boldsymbol{\theta}_0$  admits the characterization

$$\boldsymbol{\theta}_0 = \operatorname{argmax}_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} M_0(\boldsymbol{\theta}), \quad M_0(\boldsymbol{\theta}) = \lim_{n \rightarrow \infty} M_n(\boldsymbol{\theta}), \quad M_n(\boldsymbol{\theta}) = \mathbb{E}[m_n(\mathbf{z}, \boldsymbol{\theta})].$$

In other words, in the setting considered in this section,  $\hat{\boldsymbol{\theta}}_n$  approximately maximizes a function  $\hat{M}_n$  whose population counterpart  $M_n$  can be interpreted as a regularization (in the sense of [Bickel and Li, 2006](#)) of a function  $M_0$  whose maximizer  $\boldsymbol{\theta}_0$  is the object of interest. This generalization is attractive because it allows us to formulate results that cover local  $M$ -estimators such as the conditional maximum score estimator of [Honoré and Kyriazidou \(2000\)](#). Studying this setting, [Seo and Otsu \(2018\)](#) gave conditions under which  $\hat{\boldsymbol{\theta}}_n$  converges at a rate equal to the cube root of the “effective” sample size and has a limiting distribution of [Chernoff \(1964\)](#) type. Analogous conclusions will be drawn below, albeit under slightly different conditions.

For any  $n$  and any  $\delta > 0$ , define

$$\bar{m}_n(\mathbf{z}) = \sup_{m \in \mathcal{M}_n} |m(\mathbf{z})|, \quad \mathcal{M}_n = \{m_n(\cdot, \boldsymbol{\theta}) : \boldsymbol{\theta} \in \boldsymbol{\Theta}\},$$

and

$$\bar{d}_n^\delta(\mathbf{z}) = \sup_{d \in \mathcal{D}_n^\delta} |d(\mathbf{z})|, \quad \mathcal{D}_n^\delta = \{m_n(\cdot, \boldsymbol{\theta}) - m_n(\cdot, \boldsymbol{\theta}_0) : \boldsymbol{\theta} \in \boldsymbol{\Theta}_0^\delta\}, \quad \boldsymbol{\Theta}_0^\delta = \{\boldsymbol{\theta} \in \boldsymbol{\Theta} : \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \leq \delta\}.$$

**Condition CRA (Cube Root Asymptotics)** For some  $q_n > 0$  with  $r_n = \sqrt[3]{nq_n} \rightarrow \infty$ , the following are satisfied:

- (i)  $\{\mathcal{M}_n : n \geq 1\}$  is uniformly manageable for the envelopes  $\bar{m}_n$  and  $q_n \mathbb{E}[\bar{m}_n(\mathbf{z})^2] = O(1)$ .

Also,  $\sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} |M_n(\boldsymbol{\theta}) - M_0(\boldsymbol{\theta})| = o(1)$  and, for every  $\delta > 0$ ,  $\sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta} \setminus \boldsymbol{\Theta}_0^\delta} M_0(\boldsymbol{\theta}) < M_0(\boldsymbol{\theta}_0)$ .

(ii)  $\boldsymbol{\theta}_0$  is an interior point of  $\boldsymbol{\Theta}$  and, for some  $\delta > 0$ ,  $M_0$  and  $M_n$  are twice continuously differentiable on  $\boldsymbol{\Theta}_0^\delta$  and  $\sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_0^\delta} \|\partial^2[M_n(\boldsymbol{\theta}) - M_0(\boldsymbol{\theta})]/\partial\boldsymbol{\theta}\partial\boldsymbol{\theta}'\| = o(1)$ .

Also,  $r_n \|\partial M_n(\boldsymbol{\theta}_0)/\partial\boldsymbol{\theta}\| = o(1)$  and  $\mathbf{H}_0 = -\partial^2 M_0(\boldsymbol{\theta}_0)/\partial\boldsymbol{\theta}\partial\boldsymbol{\theta}'$  is positive definite.

(iii) For some  $\delta > 0$ ,  $\{\mathcal{D}_n^{\delta'} : n \geq 1, 0 < \delta' \leq \delta\}$  is uniformly manageable for the envelopes  $\bar{d}_n^{\delta'}$  and  $q_n \sup_{0 < \delta' \leq \delta} \mathbb{E}[\bar{d}_n^{\delta'}(\mathbf{z})^2/\delta'] = O(1)$ .

(iv) For every  $\delta_n > 0$  with  $\delta_n = O(r_n^{-1})$ ,  $q_n^3 r_n^{-1} \mathbb{E}[\bar{d}_n^{\delta_n}(\mathbf{z})^4] = o(1)$  and, for all  $\mathbf{s}, \mathbf{t} \in \mathbb{R}^d$  and for some  $\mathcal{C}_0$  with  $\mathcal{C}_0(\mathbf{s}, \mathbf{s}) + \mathcal{C}_0(\mathbf{t}, \mathbf{t}) - 2\mathcal{C}_0(\mathbf{s}, \mathbf{t}) > 0$  for  $\mathbf{s} \neq \mathbf{t}$ ,

$$\sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_0^{\delta_n}} \left| \frac{q_n}{\delta_n} \mathbb{E}[\{m_n(\mathbf{z}, \boldsymbol{\theta} + \delta_n \mathbf{s}) - m_n(\mathbf{z}, \boldsymbol{\theta})\} \{m_n(\mathbf{z}, \boldsymbol{\theta} + \delta_n \mathbf{t}) - m_n(\mathbf{z}, \boldsymbol{\theta})\}] - \mathcal{C}_0(\mathbf{s}, \mathbf{t}) \right| = o(1).$$

(v) For every  $\delta_n > 0$  with  $\delta_n = O(r_n^{-1})$ ,

$$\lim_{C \rightarrow \infty} \limsup_{n \rightarrow \infty} \sup_{0 < \delta \leq \delta_n} q_n \mathbb{E}[\mathbf{1}\{q_n \bar{d}_n^\delta(\mathbf{z}) > C\} \bar{d}_n^\delta(\mathbf{z})^2/\delta] = 0$$

and  $\sup_{\boldsymbol{\theta}, \boldsymbol{\theta}' \in \boldsymbol{\Theta}_0^{\delta_n}} \mathbb{E}[|m_n(\mathbf{z}, \boldsymbol{\theta}) - m_n(\mathbf{z}, \boldsymbol{\theta}')|]/\|\boldsymbol{\theta} - \boldsymbol{\theta}'\| = O(1)$ .

To interpret Condition CRA, consider first the benchmark case where  $m_n = m_0$  and  $q_n = 1$ . In this case, the condition is similar to assumptions (ii)-(vii) of the main theorem of [Kim and Pollard \(1990\)](#), to which the reader is referred for a definition of the term (uniformly) manageable. The differences between their assumptions and Condition CRA are technical in nature, since we need to slightly strengthen their assumptions in order to be able to analyze the bootstrap. For instance, the displayed part of Condition CRA(iv) is a locally uniform (with respect to  $\boldsymbol{\theta}$  near  $\boldsymbol{\theta}_0$ ) version of its counterpart in [Kim and Pollard \(1990\)](#). More generally, Condition CRA can be interpreted as an  $n$ -varying version of a suitably (for the purpose of analyzing the bootstrap) strengthened version of the assumptions of [Kim and Pollard \(1990\)](#). The differences between Condition CRA and the *i.i.d.* version of the conditions in [Seo and Otsu \(2018\)](#) are also technical in nature, but for completeness we highlight two here. First, they control the complexity of various function classes using the concept of bracketing entropy, while we follow [Kim and Pollard \(1990\)](#) and obtain maximal inequalities using bounds on uniform entropy numbers implied by the concept of (uniform) manageability. Second, whereas [Seo and Otsu \(2018\)](#) control the bias of  $\hat{\boldsymbol{\theta}}_n$  through a locally uniform bound on  $M_n - M_0$ , Condition CRA controls the bias through the first and second derivatives of  $M_n - M_0$ .

Under Condition CRA, the effective sample size is  $nq_n = r_n^3$  and if  $\hat{\boldsymbol{\theta}}_n$  is an approximate

maximizer of  $\hat{M}_n$ , then  $r_n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$  has a limiting distribution of Chernoff (1964) type. The heuristics of the previous section are rate-adaptive (i.e.,  $\sqrt[3]{n}$  can be replaced by a generic  $r_n$ ), so once again it stands to reason that if  $\tilde{\mathbf{H}}_n$  is a consistent estimator of  $\mathbf{H}_0$ , then the distribution of  $r_n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$  can be consistently estimated by that of  $r_n(\tilde{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n)$ , where  $\tilde{\boldsymbol{\theta}}_n^*$  is an approximate maximizer of

$$\tilde{M}_n^*(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \tilde{m}_n(\mathbf{z}_{i,n}^*, \boldsymbol{\theta}), \quad \tilde{m}_n(\mathbf{z}, \boldsymbol{\theta}) = m_n(\mathbf{z}, \boldsymbol{\theta}) - \hat{M}_n(\boldsymbol{\theta}) - \frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n)' \tilde{\mathbf{H}}_n (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n),$$

with  $\mathbf{z}_{1,n}^*, \dots, \mathbf{z}_{n,n}^*$  being a random sample from the empirical distribution of  $\mathbf{z}_1, \dots, \mathbf{z}_n$ . A precise statement is given in the following theorem.

**Theorem 1** *Suppose Condition CRA holds. If  $\tilde{\mathbf{H}}_n \rightarrow_{\mathbb{P}} \mathbf{H}_0$  and if*

$$\hat{M}_n(\hat{\boldsymbol{\theta}}_n) \geq \sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \hat{M}_n(\boldsymbol{\theta}) - o_{\mathbb{P}}(r_n^{-2}) \quad \text{and} \quad \tilde{M}_n^*(\tilde{\boldsymbol{\theta}}_n^*) \geq \sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \tilde{M}_n^*(\boldsymbol{\theta}) - o_{\mathbb{P}}(r_n^{-2}),$$

then

$$r_n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \rightsquigarrow \operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} \{\mathcal{G}_0(\mathbf{s}) + \mathcal{Q}_0(\mathbf{s})\}, \quad (9)$$

and

$$r_n(\tilde{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n) \rightsquigarrow_{\mathbb{P}} \operatorname{argmax}_{\mathbf{s} \in \mathbb{R}^d} \{\mathcal{G}_0(\mathbf{s}) + \mathcal{Q}_0(\mathbf{s})\}, \quad (10)$$

where  $\mathcal{G}_0$  is a zero-mean Gaussian process with covariance kernel  $\mathcal{C}_0$  and  $\mathcal{Q}_0(\mathbf{s}) = -\mathbf{s}' \mathbf{H}_0 \mathbf{s} / 2$ .

The algorithm for our proposed bootstrap-based distributional approximation is as follows:

*Step 1.* Using the sample  $\mathbf{z}_1, \dots, \mathbf{z}_n$ , compute  $\hat{\boldsymbol{\theta}}_n$  by approximately maximizing  $\hat{M}_n(\boldsymbol{\theta})$ .

*Step 2.* Using  $\hat{\boldsymbol{\theta}}_n$  and  $\mathbf{z}_1, \dots, \mathbf{z}_n$ , compute  $\tilde{\mathbf{H}}_n$ . (A generic estimator  $\tilde{\mathbf{H}}_n$  is given in Section 3.1.)

*Step 3.* Using  $\hat{\boldsymbol{\theta}}_n$ ,  $\tilde{\mathbf{H}}_n$ , and the bootstrap sample  $\mathbf{z}_{1,n}^*, \dots, \mathbf{z}_{n,n}^*$ , compute  $\tilde{\boldsymbol{\theta}}_n^*$  by approximately maximizing  $\tilde{M}_n^*(\boldsymbol{\theta})$ . ( $\hat{\boldsymbol{\theta}}_n$  and  $\tilde{\mathbf{H}}_n$  are not recomputed at this step.)

*Step 4.* Repeat *Step 3* to generate draws from the distribution of  $r_n(\tilde{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n)$ .

### 3.1 Estimation of $\mathbf{H}_0$

A generic numerical derivative estimator of  $\mathbf{H}_0$  is the matrix  $\tilde{\mathbf{H}}_n^{\text{ND}}$  with element  $(k, l)$  given by

$$\tilde{H}_{n,kl}^{\text{ND}} = -\frac{1}{4\epsilon_n^2} [\hat{M}_n(\hat{\boldsymbol{\theta}}_n + \mathbf{e}_k \epsilon_n + \mathbf{e}_l \epsilon_n) - \hat{M}_n(\hat{\boldsymbol{\theta}}_n + \mathbf{e}_k \epsilon_n - \mathbf{e}_l \epsilon_n) - \hat{M}_n(\hat{\boldsymbol{\theta}}_n - \mathbf{e}_k \epsilon_n + \mathbf{e}_l \epsilon_n) + \hat{M}_n(\hat{\boldsymbol{\theta}}_n - \mathbf{e}_k \epsilon_n - \mathbf{e}_l \epsilon_n)],$$

where  $\mathbf{e}_k$  is the  $k$ th unit vector in  $\mathbb{R}^d$  and where  $\epsilon_n$  is a positive tuning parameter. Conditions under which this estimator is consistent are given in the following lemma.

**Lemma 1** *Suppose Condition CRA holds and that  $r_n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = O_{\mathbb{P}}(1)$ . If  $\epsilon_n \rightarrow 0$  and if  $r_n \epsilon_n \rightarrow \infty$ , then  $\tilde{\mathbf{H}}_n^{\text{ND}} \rightarrow_{\mathbb{P}} \mathbf{H}_0$ .*

Plausibility of the high-level condition  $r_n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = O_{\mathbb{P}}(1)$  follows from (9). To facilitate practical implementation, it is useful to go beyond consistency and develop a Nagar-type mean squared error (MSE) expansion that can be used to select  $\epsilon_n$ . To state one such result for  $\tilde{H}_{n,kl}^{\text{ND}}$ , define

$$\begin{aligned}\ddot{M}_{n,kl}(\boldsymbol{\theta}) &= \frac{\partial^2}{\partial \theta_k \partial \theta_l} M_n(\boldsymbol{\theta}), & \ddot{M}_{0,kl}(\boldsymbol{\theta}) &= \frac{\partial^2}{\partial \theta_k \partial \theta_l} M_0(\boldsymbol{\theta}), \\ \mathbf{B}_{kl} &= -\frac{1}{6} \left[ \frac{\partial^2}{\partial \theta_k^2} \ddot{M}_{0,kl}(\boldsymbol{\theta}_0) + \frac{\partial^2}{\partial \theta_l^2} \ddot{M}_{0,kl}(\boldsymbol{\theta}_0) \right],\end{aligned}$$

and

$$\mathbf{V}_{kl} = \frac{1}{8} [\mathcal{C}_0(\mathbf{e}_k + \mathbf{e}_l, \mathbf{e}_k + \mathbf{e}_l) + \mathcal{C}_0(\mathbf{e}_k - \mathbf{e}_l, \mathbf{e}_k - \mathbf{e}_l) - 2\mathcal{C}_0(\mathbf{e}_k + \mathbf{e}_l, \mathbf{e}_k - \mathbf{e}_l) - 2\mathcal{C}_0(\mathbf{e}_k + \mathbf{e}_l, -\mathbf{e}_k + \mathbf{e}_l)].$$

**Lemma 2** *Suppose the conditions of Lemma 1 hold and that, for some  $\delta > 0$ ,  $\ddot{M}_{0,kl}$  and  $\ddot{M}_{n,kl}$  are twice continuously differentiable on  $\boldsymbol{\Theta}_0^\delta$  with  $\sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}_0^\delta} \|\partial^2[\ddot{M}_{n,kl}(\boldsymbol{\theta}) - \ddot{M}_{0,kl}(\boldsymbol{\theta})]/\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'\| = o(1)$ . If  $\mathcal{C}_0(\mathbf{s}, -\mathbf{s}) = 0$  and  $\mathcal{C}_0(\mathbf{s}, \mathbf{t}) = \mathcal{C}_0(-\mathbf{s}, -\mathbf{t})$  for all  $\mathbf{s}, \mathbf{t} \in \mathbb{R}^d$ , then  $\tilde{H}_{n,kl}^{\text{ND}}$  admits an approximation  $\check{H}_{n,kl}^{\text{ND}}$  satisfying*

$$\tilde{H}_{n,kl}^{\text{ND}} = \check{H}_{n,kl}^{\text{ND}} + o_{\mathbb{P}}\left(\epsilon_n^2 + \frac{1}{\sqrt{r_n^3 \epsilon_n^3}}\right) + O_{\mathbb{P}}\left(\frac{1}{r_n}\right),$$

where the  $O_{\mathbb{P}}(1/r_n)$  term does not depend on  $\epsilon_n$  and where

$$\mathbb{E}[(\check{H}_{n,kl}^{\text{ND}} - H_{n,kl})^2] = \epsilon_n^4 \mathbf{B}_{kl}^2 + \frac{1}{r_n^3 \epsilon_n^3} \mathbf{V}_{kl} + o\left(\epsilon_n^4 + \frac{1}{r_n^3 \epsilon_n^3}\right), \quad H_{n,kl} = -\ddot{M}_{n,kl}(\boldsymbol{\theta}_0).$$

The conditions  $\mathcal{C}_0(\mathbf{s}, -\mathbf{s}) = 0$  and  $\mathcal{C}_0(\mathbf{s}, \mathbf{t}) = \mathcal{C}_0(-\mathbf{s}, -\mathbf{t})$  are satisfied in all of the examples we have analyzed. Using the lemma, the approximate MSE (AMSE),  $\epsilon_n^4 \mathbf{B}_{kl}^2 + r_n^{-3} \epsilon_n^{-3} \mathbf{V}_{kl}$ , can be minimized by choosing  $\epsilon_n$  proportional to  $r_n^{-3/7}$ , the optimal factor of proportionality being a function of  $\mathbf{B}_{kl}$  and  $\mathbf{V}_{kl}$ . To be specific, the optimal  $\epsilon_n$  is given by  $\epsilon_{n,kl}^{\text{AMSE}} = (3\mathbf{V}_{kl}/4\mathbf{B}_{kl}^2)^{1/7} r_n^{-3/7}$ , a feasible version of which can be constructed by replacing  $\mathbf{B}_{kl}$  and  $\mathbf{V}_{kl}$  with preliminary estimators thereof.

## 4 Examples

### 4.1 Maximum Score

To describe a version of the maximum score estimator of [Manski \(1975, 1985\)](#), suppose  $\mathbf{z}_1, \dots, \mathbf{z}_n$  is a random sample of  $\mathbf{z} = (y, \mathbf{x}')'$  generated by the binary response model

$$y = \mathbb{1}(\mathbf{x}'\boldsymbol{\beta}_0 + u \geq 0), \quad \text{Median}(u|\mathbf{x}) = 0,$$

where  $\boldsymbol{\beta}_0 \in \mathbb{R}^{d+1}$  is an unknown parameter of interest,  $\mathbf{x} \in \mathbb{R}^{d+1}$  is a vector of covariates, and  $u$  is an unobserved error term. Following [Abrevaya and Huang \(2005\)](#), we employ the parameterization  $\boldsymbol{\beta}_0 = (1, \boldsymbol{\theta}_0')'$ , where  $\boldsymbol{\theta}_0 \in \mathbb{R}^d$  is unknown. In other words, we assume that the first element of  $\boldsymbol{\beta}_0$  is positive and then normalize the (unidentified) scale of  $\boldsymbol{\beta}_0$  by setting its first element equal to unity. Partitioning  $\mathbf{x}$  conformably with  $\boldsymbol{\beta}_0$  as  $\mathbf{x} = (x_1, \mathbf{x}_2')'$ , a maximum score estimator of  $\boldsymbol{\theta}_0$  is any  $\hat{\boldsymbol{\theta}}_n^{\text{MS}}$  approximately maximizing  $\hat{M}_n$  for  $m_n(\mathbf{z}, \boldsymbol{\theta}) = m^{\text{MS}}(\mathbf{z}, \boldsymbol{\theta}) = (2y - 1)\mathbb{1}(x_1 + \mathbf{x}_2'\boldsymbol{\theta} \geq 0)$ .

Regarded as a member of the class of  $M$ -estimators exhibiting cube root asymptotics, the maximum score estimator is representative in a couple of respects. First, under easy-to-interpret primitive conditions the estimator is covered by the results of [Section 3](#). Second, in addition to the generic estimator  $\tilde{\mathbf{H}}_n^{\text{ND}}$  discussed above, the maximum score estimator admits example-specific consistent estimators of the  $\mathbf{H}_0$  associated with it.

Under standard regularity conditions (stated in [Section A.2](#) of the supplemental appendix), Condition CRA is satisfied with  $q_n = 1$ ,

$$\mathbf{H}_0 = \mathbf{H}^{\text{MS}} = 2\mathbb{E}[f_{u|x_1, \mathbf{x}_2}(0 | -\mathbf{x}_2'\boldsymbol{\theta}_0, \mathbf{x}_2)f_{x_1|\mathbf{x}_2}(-\mathbf{x}_2'\boldsymbol{\theta}_0|\mathbf{x}_2)\mathbf{x}_2\mathbf{x}_2'],$$

and

$$\mathcal{C}_0(\mathbf{s}, \mathbf{t}) = \mathcal{C}^{\text{MS}}(\mathbf{s}, \mathbf{t}) = \mathbb{E}[f_{x_1|\mathbf{x}_2}(-\mathbf{x}_2'\boldsymbol{\theta}_0|\mathbf{x}_2) \min\{|\mathbf{x}_2'\mathbf{s}|, |\mathbf{x}_2'\mathbf{t}|\} \mathbb{1}(\text{sgn}(\mathbf{x}_2'\mathbf{s}) = \text{sgn}(\mathbf{x}_2'\mathbf{t}))],$$

where  $f_{a|\mathbf{b}}$  denotes the conditional Lebesgue density of  $a$  given  $\mathbf{b}$ . As a consequence, [Theorem 1](#) is applicable to  $\hat{\boldsymbol{\theta}}_n^{\text{MS}}$  and the consistency requirement  $\tilde{\mathbf{H}}_n \rightarrow_{\mathbb{P}} \mathbf{H}^{\text{MS}}$  is satisfied by the numerical derivative estimator discussed in [Section 3.1](#) if  $\epsilon_n \rightarrow 0$  and  $n\epsilon_n^3 \rightarrow \infty$ . Under the additional regularity conditions of [Lemma 2](#), MSE-optimal tuning parameter choices are feasible. In addition, alternative consistent estimators of  $\mathbf{H}^{\text{MS}}$  can be constructed exploiting the specific structure of this example. One option is to employ a “plug-in” estimator of  $\mathbf{H}^{\text{MS}}$  based on nonparametric estimators of  $f_{u|x_1, \mathbf{x}_2}$

and  $f_{x_1|x_2}$ . An alternative, example-specific estimator is

$$\tilde{\mathbf{H}}_n^{\text{MS}} = -\frac{1}{n} \sum_{i=1}^n (2y_i - 1) \dot{K}_n(x_{1i} + \mathbf{x}'_{2i} \hat{\boldsymbol{\theta}}_n^{\text{MS}}) \mathbf{x}_{2i} \mathbf{x}'_{2i},$$

where, for a differentiable kernel function  $K$  and a positive bandwidth  $h_n$ ,  $\dot{K}_n(u) = dK_n(u)/du$  and  $K_n(u) = K(u/h_n)/h_n$ . In words,  $\tilde{\mathbf{H}}_n^{\text{MS}}$  is constructed by “smoothing out” the indicator function entering  $m^{\text{MS}}(\mathbf{z}, \boldsymbol{\theta})$  and then twice differentiating the corresponding objective function (previously used by [Horowitz, 1992](#)).

## 4.2 Panel Maximum Score

Consider the panel data binary response model

$$Y_t = \mathbb{1}(\mathbf{X}'_t \boldsymbol{\beta}_0 + \alpha + u_t \geq 0), \quad t = 1, 2,$$

where  $\boldsymbol{\beta}_0 \in \mathbb{R}^{d+1}$  is an unknown parameter of interest,  $\alpha$  is an unobserved (time-invariant) individual-specific effect, and  $u_t$  is an unobserved error term. Analyzing this model, [Manski \(1987\)](#) gave conditions under which  $\boldsymbol{\beta}_0$  is identified up to scale and demonstrated consistency of a conditional maximum score estimator.

Suppose  $\boldsymbol{\beta}_0$  is identified up to scale and that its first element is positive, in which case we can normalize that element to unity and employ the parameterization  $\boldsymbol{\beta}_0 = (1, \boldsymbol{\theta}'_0)'$ , where  $\boldsymbol{\theta}_0 \in \mathbb{R}^d$  is unknown. To describe a version of the estimator of [Manski \(1987\)](#), partition  $\mathbf{X}_t$  conformably with  $\boldsymbol{\beta}_0$  as  $\mathbf{X}_t = (X_{1t}, \mathbf{X}'_{2t})'$  and define  $\mathbf{z} = (y, x_1, \mathbf{x}'_2)'$ , where  $y = Y_2 - Y_1$ ,  $x_1 = X_{12} - X_{11}$ , and  $\mathbf{x}_2 = (\mathbf{X}_{22} - \mathbf{X}_{21})$ . Assuming  $\mathbf{z}_1, \dots, \mathbf{z}_n$  is a random sample of  $\mathbf{z}$ , a panel maximum score estimator of  $\boldsymbol{\theta}_0$  is any  $\hat{\boldsymbol{\theta}}_n^{\text{PMS}}$  approximately maximizing  $\hat{M}_n$  for  $m_n(\mathbf{z}, \boldsymbol{\theta}) = m^{\text{PMS}}(\mathbf{z}, \boldsymbol{\theta}) = y \mathbb{1}(x_1 + \mathbf{x}'_2 \boldsymbol{\theta} \geq 0)$ .

As one would expect, the properties of  $\hat{\boldsymbol{\theta}}_n^{\text{PMS}}$  are qualitatively similar to those of  $\hat{\boldsymbol{\theta}}_n^{\text{MS}}$ . To be specific, under regularity conditions (stated in Section A.3 of the supplemental appendix), the panel maximum score estimator is covered by the results of Section 3 and an example-specific alternative to the generic numerical derivative estimator is available, namely

$$\tilde{\mathbf{H}}_n^{\text{PMS}} = -n^{-1} \sum_{i=1}^n y_i \dot{K}_n(x_{1i} + \mathbf{x}'_{2i} \hat{\boldsymbol{\theta}}_n^{\text{PMS}}) \mathbf{x}_{2i} \mathbf{x}'_{2i},$$

where  $\dot{K}_n$  is as in the maximum score example.

### 4.3 Conditional Maximum Score

Consider the dynamic panel data binary response model

$$Y_t = \mathbb{1}(\mathbf{X}_t' \boldsymbol{\beta}_0 + Y_{t-1} \gamma_0 + \alpha + u_t \geq 0), \quad t = 1, 2, 3,$$

where  $\boldsymbol{\beta}_0 \in \mathbb{R}^d$  and  $\gamma_0 \in \mathbb{R}$  are unknown parameters of interest,  $\alpha$  is an unobserved (time-invariant) individual-specific effect, and  $u_t$  is an unobserved error term. [Honoré and Kyriazidou \(2000\)](#) analyzed this model and gave conditions under which  $\boldsymbol{\beta}_0$  and  $\gamma_0$  are identified up to a common scale factor. Assuming these conditions hold and that the first element of  $\boldsymbol{\beta}_0$  is positive, we can normalize that element to unity and employ the parameterization  $(\boldsymbol{\beta}_0', \gamma_0)' = (1, \boldsymbol{\theta}_0)'$ , where  $\boldsymbol{\theta}_0 \in \mathbb{R}^d$  is unknown.

To describe a version of the conditional maximum score estimator of [Honoré and Kyriazidou \(2000\)](#), partition  $\mathbf{X}_t$  after the first element as  $\mathbf{X}_t = (X_{1t}, \mathbf{X}_{2t}')'$  and define  $\mathbf{z} = (y, x_1, \mathbf{x}_2', \mathbf{w}')'$ , where  $y = Y_2 - Y_1$ ,  $x_1 = X_{12} - X_{11}$ ,  $\mathbf{x}_2 = ((\mathbf{X}_{22} - \mathbf{X}_{21})', Y_3 - Y_0)'$ , and  $\mathbf{w} = \mathbf{X}_2 - \mathbf{X}_3$ . Assuming  $\mathbf{z}_1, \dots, \mathbf{z}_n$  is a random sample of  $\mathbf{z}$ , a conditional maximum score estimator of  $\boldsymbol{\theta}_0$  is any  $\hat{\boldsymbol{\theta}}_n^{\text{CMS}}$  approximately maximizing  $\hat{M}_n$  for  $m_n(\mathbf{z}, \boldsymbol{\theta}) = m_n^{\text{CMS}}(\mathbf{z}, \boldsymbol{\theta}) = y \mathbb{1}(x_1 + \mathbf{x}_2' \boldsymbol{\theta} \geq 0) \kappa_n(\mathbf{w})$ , where  $\kappa_n(\mathbf{w}) = \kappa(\mathbf{w}/b_n)/b_n^d$  for a kernel function  $\kappa$  and a bandwidth  $b_n$ .

Through its dependence on  $b_n$ , the function  $m_n^{\text{CMS}}$  depends on  $n$  in a non-negligible way. In particular, the effective sample size is  $nb_n^d$  (rather than  $n$ ) in the current setting, so to the extent that they exist one would expect primitive sufficient conditions for Condition CRA to include  $q_n = b_n^d$  in this example. Apart from this predictable change, the properties of the conditional maximum score estimator  $\hat{\boldsymbol{\theta}}_n^{\text{CMS}}$  turn out to be qualitatively similar to those of  $\hat{\boldsymbol{\theta}}_n^{\text{MS}}$ . To be specific, under regularity conditions (stated in Section A.4 of the supplemental appendix), the conditional maximum score estimator is covered by the results of Section 3 and an example-specific alternative to the generic numerical derivative estimator is available, namely

$$\tilde{\mathbf{H}}_n^{\text{CMS}} = -n^{-1} \sum_{i=1}^n y_i \dot{K}_n(x_{1i} + \mathbf{x}_{2i}' \hat{\boldsymbol{\theta}}_n^{\text{CMS}}) \mathbf{x}_{2i} \mathbf{x}_{2i}' \kappa_n(\mathbf{w}_i),$$

where  $\dot{K}_n$  is as in the maximum score example.

### 4.4 Empirical Risk Minimization

[Mohammadi and van de Geer \(2005\)](#) considered two-category classification problems in machine

learning. Specifically, given a binary outcome  $y \in \{-1, 1\}$  and a vector of features  $\mathbf{x} \in \mathcal{X}$ , the goal is to estimate the  $\boldsymbol{\theta}_0$  that minimizes the misclassification error (or risk)  $\mathbb{P}[h_{\boldsymbol{\theta}}(\mathbf{x}) \neq y]$  with respect to  $\boldsymbol{\theta} \in \boldsymbol{\Theta} \subseteq \mathbb{R}^d$ , where  $\{h_{\boldsymbol{\theta}} : \boldsymbol{\theta} \in \boldsymbol{\Theta}\}$  is a collection of classifiers. For simplicity, we consider the case where the feature is univariate with support  $\mathcal{X} = [0, 1]$  and the classifiers are of the form

$$h_{\boldsymbol{\theta}}(x) = \sum_{\ell=1}^{d+1} (-1)^{\ell} \mathbf{1}(\theta_{\ell-1} \leq x < \theta_{\ell}), \quad \boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_d)',$$

where  $\boldsymbol{\Theta} = \{(\theta_1, \theta_2, \dots, \theta_d)' \in [0, 1]^d : 0 = \theta_0 \leq \theta_1 \leq \dots \leq \theta_d \leq \theta_{d+1} = 1\}$ .

Assuming  $\mathbf{z}_1, \dots, \mathbf{z}_n$  is a random sample of  $\mathbf{z}$ , an empirical risk minimizer is any  $\hat{\boldsymbol{\theta}}_n^{\text{ERM}}$  approximately maximizing  $\hat{M}_n$  for  $m_n(\mathbf{z}, \boldsymbol{\theta}) = m^{\text{ERM}}(\mathbf{z}, \boldsymbol{\theta}) = -\mathbf{1}(h_{\boldsymbol{\theta}}(x) \neq y)$ . Under regularity conditions similar to those of [Mohammadi and van de Geer \(2005, Section 2.1\)](#), the empirical risk minimizer is covered by Theorem 1 and the consistency requirement on  $\tilde{\mathbf{H}}_n$  can be met in various ways; for details, see Section A.5 of the supplemental appendix.

## 5 Simulations

We illustrate the numerical performance of our proposed bootstrap-based inference methods for the maximum score estimator. Given the setup in Section 4.1, we generate data from that model with  $d = 1$ ,  $\boldsymbol{\theta}_0 = 1$ ,  $\mathbf{x} = (x_1, x_2)' \sim \mathcal{N}((0, 1)', \mathbf{I}_2)$  with  $\mathbf{I}_2$  the  $(2 \times 2)$  identity matrix, and  $u$  generated by three distinct distributions. Specifically, DGP 1 sets  $u \sim \text{Logistic}(0, 1)/\sqrt{2\pi^2/3}$ , DGP 2 sets  $u \sim \mathcal{T}_3/\sqrt{3}$ , where  $\mathcal{T}_3$  denotes a Student's  $t$ -distribution with 3 degrees of freedom, and DGP 3 sets  $u \sim (1 + 2(x_1 + x_2)^2 + (x_1 + x_2)^4)\text{Logistic}(0, 1)/\sqrt{\pi^2/48}$ .

The Monte Carlo experiment employs a sample size  $n = 1,000$  with  $B = 2,000$  bootstrap replications and  $S = 2,000$  simulations. For each of the three DGPs, we implement the standard non-parametric bootstrap, the  $m$ -out-of- $n$  bootstrap using  $m \in \{\lceil n^{1/2} \rceil, \lceil n^{2/3} \rceil, \lceil n^{4/5} \rceil\}$ , and our proposed method using the two estimators  $\tilde{\mathbf{H}}_n^{\text{MS}}$  and  $\tilde{\mathbf{H}}_n^{\text{ND}}$  of  $\mathbf{H}_0$ . We report empirical coverage for nominal 95% confidence intervals and their average interval length. For the case of our proposed procedures, we investigate their performance using (i) infeasible (simulation-based) MSE-optimal choices of tuning parameters (bandwidth/derivative step), denoted by  $h_{\text{MSE}}$  and  $\epsilon_{\text{MSE}}$ , and (ii) infeasible and feasible AMSE-optimal choices of the tuning parameters, denoted by  $h_{\text{AMSE}}$ ,  $\hat{h}_{\text{AMSE}}$ ,  $\epsilon_{\text{AMSE}}$  and  $\hat{\epsilon}_{\text{AMSE}}$ ; for details, see Section A.2 of the supplemental appendix.



Table 1: Simulations, Maximum Score Estimator, 95% Confidence Intervals.

	DGP 1			DGP 2			DGP 3		
	$h, \epsilon$	Coverage	Length	$h, \epsilon$	Coverage	Length	$h, \epsilon$	Coverage	Length
<b>Standard</b>		0.625	0.472		0.647	0.475		0.654	0.243
<b>m-out-of-n</b>									
$m = \lceil n^{1/2} \rceil$		0.997	1.698		0.998	1.753		1.000	1.890
$m = \lceil n^{2/3} \rceil$		0.978	1.185		0.983	1.221		0.989	0.724
$m = \lceil n^{4/5} \rceil$		0.899	0.820		0.897	0.837		0.930	0.447
<b>Plug-in: <math>\tilde{\mathbf{V}}_n^{\text{MS}}</math></b>									
$h_{\text{MSE}}$	0.620	0.954	0.511	0.580	0.957	0.523	0.150	0.962	0.277
$h_{\text{AMSE}}$	1.108	0.972	0.590	0.480	0.951	0.518	0.123	0.942	0.263
$\hat{h}_{\text{AMSE}}$	0.443	0.940	0.508	0.409	0.946	0.518	0.155	0.957	0.278
<b>Num Deriv: <math>\tilde{\mathbf{V}}_n^{\text{ND}}</math></b>									
$\epsilon_{\text{MSE}}$	1.400	0.936	0.483	1.360	0.938	0.485	0.290	0.939	0.249
$\epsilon_{\text{AMSE}}$	0.537	0.880	0.414	0.573	0.894	0.426	0.224	0.902	0.227
$\hat{\epsilon}_{\text{AMSE}}$	0.518	0.876	0.413	0.512	0.882	0.420	0.369	0.947	0.270

Notes:

(i) Panel **Standard** refers to standard nonparametric bootstrap, Panel **m-out-of-n** refers to  $m$ -out-of- $n$  nonparametric bootstrap with subsample size  $m$ , Panel **Plug-in:  $\tilde{\mathbf{V}}_n^{\text{MS}}$**  refers to our proposed bootstrap-based implemented using the example-specific plug-in drift estimator, and Panel **Num Deriv:  $\tilde{\mathbf{V}}_n^{\text{ND}}$**  refers to our proposed bootstrap-based implemented using the generic numerical derivative drift estimator.

(ii) Column “ $h, \epsilon$ ” reports tuning parameter value used or average across simulations when estimated, and Columns “Coverage” and “Length” report empirical coverage and average length of bootstrap-based 95% percentile confidence intervals, respectively.

(iii)  $h_{\text{MSE}}$  and  $\epsilon_{\text{MSE}}$  correspond to the simulation MSE-optimal choices,  $h_{\text{AMSE}}$  and  $\epsilon_{\text{AMSE}}$  correspond to the AMSE-optimal choices, and  $\hat{h}_{\text{AMSE}}$  and  $\hat{\epsilon}_{\text{AMSE}}$  correspond to the ROT feasible implementation of  $\hat{h}_{\text{AMSE}}$  and  $\hat{\epsilon}_{\text{AMSE}}$  described in the supplemental appendix.

Table 1 presents the main results, which are consistent across all three simulation designs. First, as expected, the standard nonparametric bootstrap (labeled “Standard”) does not perform well, leading to confidence intervals with an average 64% empirical coverage rate. Second, the  $m$ -out-of- $n$  bootstrap (labeled “m-out-of-n”) performs somewhat better for small subsamples, but leads to very large average interval length of the resulting confidence intervals. Our proposed methods, on the other hand, exhibit good finite sample performance in this Monte Carlo experiment. Results employing the example-specific plug-in estimator  $\tilde{\mathbf{H}}_n^{\text{MS}}$  are presented under the label “Plug-in” while results employing the generic numerical derivative estimator  $\tilde{\mathbf{H}}_n^{\text{ND}}$  are reported under the label “Num Deriv”. Empirical coverage appears stable across different values of the tuning parameters for each method, with better performance in the case of  $\tilde{\mathbf{H}}_n^{\text{MS}}$ . We conjecture that  $n = 1,000$  is too small for the numerical derivative estimator  $\tilde{\mathbf{H}}_n^{\text{ND}}$  to lead to as good inference performance as  $\tilde{\mathbf{H}}_n^{\text{MS}}$  (e.g., note that the MSE-optimal choice  $\epsilon_{\text{MSE}}$  is greater than 1). Nevertheless, empirical coverage of confidence

intervals constructed using our proposed bootstrap-based method is close to 95% in all cases except when  $\tilde{\mathbf{H}}_n^{\text{ND}}$  is used with either the infeasible asymptotic choice  $\epsilon_{\text{AMSE}}$  or its estimated counterpart  $\hat{\epsilon}_{\text{AMSE}}$ , and with an average interval length of at most half that of any of the  $m$ -out-of- $n$  competing confidence intervals. In particular, confidence intervals based on  $\tilde{\mathbf{H}}_n^{\text{MS}}$  implemented with the feasible bandwidth  $\hat{h}_{\text{AMSE}}$  perform quite well across the three DGPs considered.

## 6 Conclusion

We developed a valid resampling procedure for cube root asymptotics based on the nonparametric bootstrap. Whereas the standard nonparametric bootstrap is known to be invalid in the setting we study, we show that bootstrap validity can be restored by applying a carefully tailored reshaping of the objective function defining the estimator. Such reshaping is easy to implement both in general and in specific cases, as illustrated by the distinct examples we considered.

Seo and Otsu (2018) gave conditions under which results of the form (9) can be obtained also when the data exhibits weak dependence; see also Bagchi, Banerjee, and Stoev (2016), and references therein. It seems plausible that a version of our procedure, implemented with a resampling procedure suitable for dependent data, can be shown to be consistent under similar conditions, but it is beyond the scope of this paper to substantiate that conjecture.

## 7 Proof of Theorem 1

The proof proceeds by first showing (9) and then using that result to establish (10). In both cases, we employ arguments similar to those used in the proof of the main theorem of Kim and Pollard (1990). The remainder of this section outlines the main steps in the proof; for technical details, see Lemmas A.1-A.10 in Section A.1 of the supplemental appendix.

*Proof of (9).* The estimator  $\hat{\boldsymbol{\theta}}_n$  is assumed to satisfy

$$\left\{ \hat{G}_n(\mathbf{s}) + Q_n(\mathbf{s}) \right\} \Big|_{\mathbf{s} = r_n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)} \geq \sup_{\mathbf{s} \in \mathbb{R}^d} \{ \hat{G}_n(\mathbf{s}) + Q_n(\mathbf{s}) \} + o_{\mathbb{P}}(1),$$

where

$$\hat{G}_n(\mathbf{s}) = r_n^2[\hat{M}_n(\boldsymbol{\theta}_0 + \mathbf{s}r_n^{-1}) - \hat{M}_n(\boldsymbol{\theta}_0) - M_n(\boldsymbol{\theta}_0 + \mathbf{s}r_n^{-1}) + M_n(\boldsymbol{\theta}_0)]\mathbf{1}(\boldsymbol{\theta}_0 + \mathbf{s}r_n^{-1} \in \boldsymbol{\Theta})$$

and

$$Q_n(\mathbf{s}) = r_n^2[M_n(\boldsymbol{\theta}_0 + \mathbf{s}r_n^{-1}) - M_n(\boldsymbol{\theta}_0)]\mathbf{1}(\boldsymbol{\theta}_0 + \mathbf{s}r_n^{-1} \in \boldsymbol{\Theta}).$$

By the argmax continuous mapping theorem (e.g., [van der Vaart and Wellner, 1996](#), Theorem 3.2.2), it therefore suffices to show that  $r_n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = O_{\mathbb{P}}(1)$  and that  $\hat{G}_n + Q_n \rightsquigarrow \mathcal{G}_0 + \mathcal{Q}_0$  in the topology of uniform convergence on compacta. (The other conditions required by the argmax continuous mapping theorem are easily verified.)

To obtain the rate of convergence of  $\hat{\boldsymbol{\theta}}_n$ , we begin by using a standard argument to show that  $\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0 = o_{\mathbb{P}}(1)$  under Condition CRA(i) and then strengthen that conclusion to  $r_n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = O_{\mathbb{P}}(1)$  by using Conditions CRA(ii)-(iii) and proceeding along the lines of [van der Vaart and Wellner \(1996, Theorem 3.2.5\)](#). In both cases, we employ the maximal inequality in [Pollard \(1989, Theorem 4.2\)](#); for details, see Lemmas A.1 and A.3 of the supplemental appendix.

Next, because  $Q_n$  is non-random,  $\hat{G}_n + Q_n \rightsquigarrow \mathcal{G}_0 + \mathcal{Q}_0$  in the topology of uniform convergence on compacta if  $Q_n$  converges compactly to  $\mathcal{Q}_0$  and if  $\hat{G}_n \rightsquigarrow \mathcal{G}_0$  in the topology of uniform convergence on compacta. Compact convergence of  $Q_n$  follows from Condition CRA (ii); for details, see Lemma A.2 of the supplemental appendix. Also, to show that  $\hat{G}_n \rightsquigarrow \mathcal{G}_0$  in the topology of uniform convergence on compacta, it suffices to show that  $\hat{G}_n$  converges to  $\mathcal{G}_0$  in the sense of weak convergence of finite-dimensional projections and that  $\{\hat{G}_n(\mathbf{s}) : \|\mathbf{s}\| \leq K\}$  is stochastically equicontinuous for every  $K > 0$ .

Under Conditions CRA(ii)-(iv), weak convergence of finite-dimensional projections can be shown using the Cramér-Wold device and the fact that  $\mathbb{E}[\hat{G}_n(\mathbf{s})\hat{G}_n(\mathbf{t})]$  converges to  $\mathcal{C}_0(\mathbf{s}, \mathbf{t})$  for every  $\mathbf{s}, \mathbf{t} \in \mathbb{R}^d$ ; for details, see Lemma A.4 of the supplemental appendix. Finally, under Conditions CRA(iii) and CRA(v) and employing the maximal inequality in [Pollard \(1989, Theorem 4.2\)](#), stochastic equicontinuity can be shown by proceeding as in the proof of [Kim and Pollard \(1990, Lemma 4.6\)](#); for details, see Lemma A.5 of the supplemental appendix.

*Proof of (10).* The proof of (10) is a natural bootstrap analog of the proof of (9). The estimator

$\tilde{\boldsymbol{\theta}}_n^*$  is assumed to satisfy

$$\{\tilde{G}_n^*(\mathbf{s}) + \tilde{Q}_n(\mathbf{s})\} \Big|_{\mathbf{s}=r_n(\tilde{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n)} \geq \sup_{\mathbf{s} \in \mathbb{R}^d} \{\tilde{G}_n^*(\mathbf{s}) + \tilde{Q}_n(\mathbf{s})\} + o_{\mathbb{P}}(1),$$

where

$$\tilde{G}_n^*(\mathbf{s}) = r_n^2 [\tilde{M}_n^*(\hat{\boldsymbol{\theta}}_n + \mathbf{s}r_n^{-1}) - \tilde{M}_n^*(\hat{\boldsymbol{\theta}}_n) - \tilde{M}_n(\hat{\boldsymbol{\theta}}_n + \mathbf{s}r_n^{-1}) + \tilde{M}_n(\hat{\boldsymbol{\theta}}_n)] \mathbf{1}(\hat{\boldsymbol{\theta}}_n + \mathbf{s}r_n^{-1} \in \boldsymbol{\Theta})$$

and

$$\tilde{Q}_n(\mathbf{s}) = r_n^2 [\tilde{M}_n(\hat{\boldsymbol{\theta}}_n + \mathbf{s}r_n^{-1}) - \tilde{M}_n(\hat{\boldsymbol{\theta}}_n)] \mathbf{1}(\hat{\boldsymbol{\theta}}_n + \mathbf{s}r_n^{-1} \in \boldsymbol{\Theta}) = -\frac{1}{2} \mathbf{s}' \tilde{\mathbf{H}}_n \mathbf{s} \mathbf{1}(\hat{\boldsymbol{\theta}}_n + \mathbf{s}r_n^{-1} \in \boldsymbol{\Theta}).$$

By the argmax continuous mapping theorem, it therefore suffices to show that  $r_n(\tilde{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n) = O_{\mathbb{P}}(1)$  and that  $\tilde{G}_n^* + \tilde{Q}_n \rightsquigarrow_{\mathbb{P}} \mathcal{G}_0 + \mathcal{Q}_0$  in the topology of uniform convergence on compacta.

Using  $\tilde{\mathbf{H}}_n \rightarrow_{\mathbb{P}} \mathbf{H}_0$ , to obtain the rate of convergence of  $\tilde{\boldsymbol{\theta}}_n^*$  we first show that  $\tilde{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n = o_{\mathbb{P}}(1)$  under Condition CRA(i) and then strengthen that conclusion to  $r_n(\tilde{\boldsymbol{\theta}}_n^* - \hat{\boldsymbol{\theta}}_n) = O_{\mathbb{P}}(1)$  by using  $r_n(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = O_{\mathbb{P}}(1)$  and Condition CRA(iii). As in the derivation of the convergence rate of  $\hat{\boldsymbol{\theta}}_n$ , both steps employ the maximal inequality in Pollard (1989, Theorem 4.2); for details, see Lemmas A.6 and A.8 of the supplemental appendix.

Next, because  $\mathcal{Q}_0$  is non-random,  $\tilde{G}_n^* + \tilde{Q}_n \rightsquigarrow_{\mathbb{P}} \mathcal{G}_0 + \mathcal{Q}_0$  in the topology of uniform convergence on compacta if  $\tilde{Q}_n \rightarrow_{\mathbb{P}} \mathcal{Q}_0$  in the topology of uniform convergence on compacta and if  $\tilde{G}_n^* \rightsquigarrow \mathcal{G}_0$  in the topology of uniform convergence on compacta. By construction,  $\tilde{Q}_n$  is such that if  $\tilde{\mathbf{H}}_n \rightarrow_{\mathbb{P}} \mathbf{H}_0$  and if  $\hat{\boldsymbol{\theta}}_n \rightarrow_{\mathbb{P}} \boldsymbol{\theta}_0 \in \text{int}(\boldsymbol{\Theta})$ , then  $\tilde{Q}_n \rightarrow_{\mathbb{P}} \mathcal{Q}_0$  in the topology of uniform convergence on compacta; for details, see Lemma A.7 of the supplemental appendix.

Also, to show that  $\tilde{G}_n^* \rightsquigarrow_{\mathbb{P}} \mathcal{G}_0$  in the topology of uniform convergence on compacta, it suffices to show that  $\tilde{G}_n^*$  converges to  $\mathcal{G}_0$  in the sense of conditional weak convergence in probability of finite-dimensional projections and that  $\{\tilde{G}_n^*(\mathbf{s}) : \|\mathbf{s}\| \leq K\}$  is stochastically equicontinuous for every  $K > 0$ . Conditional weak convergence in probability of finite-dimensional projections can be shown using the Cramér-Wold device and the fact that the maximal inequality in Pollard (1989, Theorem 4.2) can be used to show that  $\mathbb{E}_n^*[\tilde{G}_n^*(\mathbf{s})\tilde{G}_n^*(\mathbf{t})]$  converges in probability to  $\mathcal{C}_0(\mathbf{s}, \mathbf{t})$  for every  $\mathbf{s}, \mathbf{t} \in \mathbb{R}^d$ , where  $\mathbb{E}_n^*$  denotes an expectation computed under the bootstrap distribution conditional

on the data; for details, see Lemma A.9 of the supplemental appendix. Finally, employing the maximal inequality in Pollard (1989, Theorem 4.2), stochastic equicontinuity can be shown by proceeding as in the proof of Kim and Pollard (1990, Lemma 4.6); for details, see Lemma A.10 of the supplemental appendix.

## References

- ABREVAYA, J., AND J. HUANG (2005): “On the Bootstrap of the Maximum Score Estimator,” *Econometrica*, 73(4), 1175–1204.
- ANDREWS, D. W. K., AND G. SOARES (2010): “Inference for Parameters Defined by Moment Inequalities Using Generalized Moment Selection,” *Econometrica*, 78(1), 119–157.
- BAGCHI, P., M. BANERJEE, AND S. A. STOEV (2016): “Inference for Monotone Functions under Short-and Long-range Dependence: Confidence Intervals and New Universal Limits,” *Journal of the American Statistical Association*, 111(516), 1634–1647.
- BICKEL, P. J., F. GÖTZE, AND W. R. VAN ZWET (1997): “Resampling Fewer than  $n$  Observations: Gains, Losses, and Remedies for Losses,” *Statistica Sinica*, 7(1), 1–31.
- BICKEL, P. J., AND B. LI (2006): “Regularization in Statistics,” *Test*, 15(2), 271–344.
- CHERNOFF, H. (1964): “Estimation of the Mode,” *Annals of the Institute of Statistical Mathematics*, 16(1), 31–41.
- DELGADO, M. A., J. M. RODRIGUEZ-POO, AND M. WOLF (2001): “Subsampling Inference in Cube Root Asymptotics with an Application to Manski’s Maximum Score Estimator,” *Economics Letters*, 73(2), 241–250.
- DÜMBGEN, L. (1993): “On Nondifferentiable Functions and the Bootstrap,” *Probability Theory and Related Fields*, 95(1), 125–140.
- FANG, Z., AND A. SANTOS (2019): “Inference on Directionally Differentiable Functions,” *Review of Economic Studies*, 86(1), 377–412.
- GRENANDER, U. (1956): “On the Theory of Mortality Measurement: Part II,” *Scandinavian Actuarial Journal*, 39(2), 125–153.
- GROENEBOOM, P., AND K. HENDRICKX (2018): “Current Status Linear Regression,” *Annals of Statistics*, 46(4), 1415–1444.
- GROENEBOOM, P., AND G. JONGBLOED (2018): “Some Developments in the Theory of Shape Constrained Inference,” *Statistical Science*, 33(4), 473–492.
- HONG, H., AND J. LI (2020): “The Numerical Bootstrap,” *Annals of Statistics*, 48(1), 397–412.
- HONORÉ, B. E., AND E. KYRIAZIDOU (2000): “Panel Data Discrete Choice Models with Lagged Dependent Variables,” *Econometrica*, 68(4), 839–874.
- HOROWITZ, J. L. (1992): “A Smoothed Maximum Score Estimator for the Binary Response Model,” *Econometrica*, 60(3), 505–531.

- KIM, J., AND D. POLLARD (1990): “Cube Root Asymptotics,” *Annals of Statistics*, 18(1), 191–219.
- KOSOROK, M. R. (2008): “Bootstrapping the Grenander Estimator,” in *Beyond Parametrics in Interdisciplinary Research: Festschrift in Honor of Professor Pranab K. Sen*, pp. 282–292. Institute of Mathematical Statistics.
- LEE, S. M. S., AND M. C. PUN (2006): “On  $m$  out of  $n$  Bootstrapping for Nonstandard M-Estimation With Nuisance Parameters,” *Journal of the American Statistical Association*, 101(475), 1185–1197.
- LEE, S. M. S., AND P. YANG (2020): “Bootstrap Confidence Regions Based on M-Estimators under Nonstandard Conditions,” *Annals of Statistics*, 48(1), 274–299.
- LÉGER, C., AND B. MACGIBBON (2006): “On the Bootstrap in Cube Root Asymptotics,” *Canadian Journal of Statistics*, 34(1), 29–44.
- MANSKI, C. F. (1975): “Maximum Score Estimation of the Stochastic Utility Model of Choice,” *Journal of Econometrics*, 3(3), 205–228.
- (1985): “Semiparametric Analysis of Discrete Response: Asymptotic Properties of the Maximum Score Estimator,” *Journal of Econometrics*, 27(3), 313–333.
- (1987): “Semiparametric Analysis of Random Effects Linear Models from Binary Panel Data,” *Econometrica*, 55(2), 357–362.
- MOHAMMADI, L., AND S. VAN DE GEER (2005): “Asymptotics in Empirical Risk Minimization,” *Journal of Machine Learning Research*, 6(Dec), 2027–2047.
- MUKHERJEE, D., M. BANERJEE, AND Y. RITOV (2019): “Non-Standard Asymptotics in High Dimensions: Manski’s Maximum Score Estimator Revisited,” arXiv:1903.10063.
- PATRA, R. K., E. SEIJO, AND B. SEN (2018): “A Consistent Bootstrap Procedure for the Maximum Score Estimator,” *Journal of Econometrics*, 205(2), 488–507.
- POLITIS, D. N., AND J. P. ROMANO (1994): “Large Sample Confidence Regions Based on Subsamples under Minimal Assumptions,” *Annals of Statistics*, 22(4), 2031–2050.
- POLLARD, D. (1989): “Asymptotics via Empirical Processes,” *Statistical Science*, 4(4), 341–354.
- SEN, B., M. BANERJEE, AND M. WOODROOFE (2010): “Inconsistency of Bootstrap: The Grenander Estimator,” *Annals of Statistics*, 38(4), 1953–1977.
- SEO, M. H., AND T. OTSU (2018): “Local M-Estimation with Discontinuous Criterion for Dependent and Limited Observations,” *Annals of Statistics*, 46(1), 344–369.
- SHI, C., W. LU, AND R. SONG (2018): “A Massive Data Framework for M-estimators with Cubic Rate,” *Journal of the American Statistical Association*, 113(524), 1698–1709.
- VAN DER VAART, A. W., AND J. A. WELLNER (1996): *Weak Convergence and Empirical Processes*. Springer.